# Document Forensics as a Unique IDP Component

## Build your document processing flow with expert components
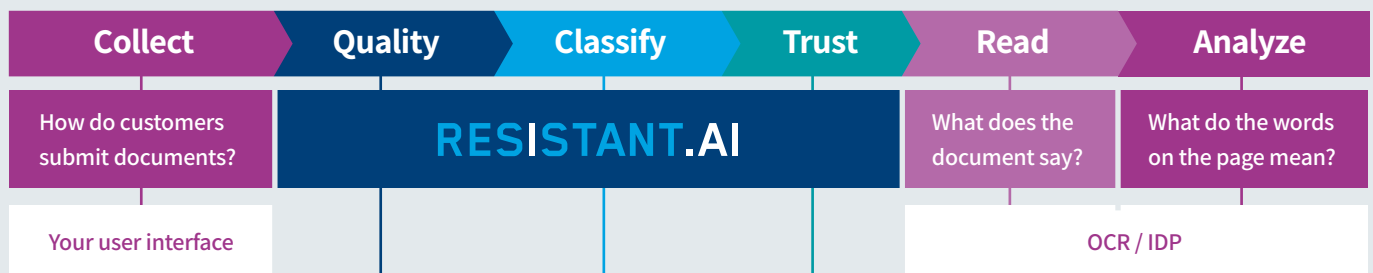
You're growing.

Maybe you're KYC onboarding in a new geographic market, or maybe you're starting to offer services with new compliance requirements. No matter the challenges your company is tackling, you need smooth document processing to keep moving up. But is the optical character recognition (OCR) solution that converts your customers' documents to information you can use equally fluent in English and Mandarin? And does your intelligent document processor (IDP) handle ID cards and invoices equally well? More than likely, document processing that is painless for you and for your customers requires a mix of technical solutions with unique areas of expertise.

Our area of expertise is fraud: an artificial intelligence layer that augments your existing flow to ensure the documents you're receiving are both usable and legitimate before you incur costs—and potential risks—down the line. Remove frustration for genuine customers and stop fraud at the source with Document Forensics.

## First fight fraud

You handle document collection and document processors get you the info you need. We sit in between, reducing OCR costs on unreadable or irrelevant documents—and acting as your first line of defense against bad actors.

| Collect | Quality | Classify | Trust | Read | Analyze |
|---------|---------|----------|-------|------|---------|
| How do customers submit documents? | | RESISTANT.AI | | What does the document say? | What do the words on the page mean? |
| Your user interface | | | | OCR / IDP | |

### Quality:

**Is this document usable?**

Research confirms that accurate OCR requires high-quality images,[1] so sending unreadable documents downstream only incurs costs with no benefits. Worse still, poor image quality is a favorite technique for hiding manipulated and outright forged documents.

Resistant AI checks for flash reflections, blur, low resolution, and more within three seconds, rejecting unusable or suspicious documents.

Feedback can then immediately prompt customers to re-upload new documents, reducing frustrating and futile wait times while adding friction for fraudsters.

### Classify:

**Is this the right type of document?**

With as few as 50 to 80 samples, our AI learns what a given document from a given issuer should look like in both images and PDFs. This lets you leverage an adaptable knowledge base to set automatic acceptance criteria for a nearly endless number of document types—accept from this issuer but not that one, double-check that this entity actually exists, accept photos but not scans, etc.

This principle can also help streamline and augment the rest of your documents flow, from choosing the right OCR/IDP provider to highlighting "hotspots" of relevant content.

### Trust:

**Is this document genuine?**

Our AI anomaly detection evaluates over 500 characteristics applicable across document types—and often invisible to the human eye.

We verify that an appropriate author, program, and device produced a document. New documents are compared to historical ones to flag reused documents and counterfeiting templates—often calling cards of serial fraud rings. Features including logo detection confirm customer documents match how known institutions create legitimate documents. And that's just the start.

Verified documents can be automatically approved, again saving unnecessary processing and post-hoc fraud checks.

[1] Karez Hamad, et al. "A Detailed Analysis of Optical Character Recognition Technology" (International Journal of Applied Mathematics, Electronics and Computers, 2016)

## Structure and specialization: why choosing an all-in-one IDP is hard

With data capture accuracy averaging 95%,[2] today's technologies fare well when supplied with high-quality documents. But OCR and more complex IDPs need to navigate two intertwined and compounding obstacles to achieve their fullest potential: localization—understanding the language a document is written in—and parsing—recognizing information on the page that is actually important and organizing it in a useful way.

The table below shows the results of a study on text extraction accuracy.[3] It's far from comprehensive, but the findings are applicable across the market: different solutions have different strengths and weaknesses based on quality and document type. Success with handwritten and multilingual documents are particularly variable.

Inaccuracies in text extracted from images can of course be magnified downstream, but parsing is also a major issue even with native PDFs, where content is extracted nearly perfectly at the code level. Structured data such as government-issued forms, semi-structured data such as common invoices, and unstructured data such as one-off letters all challenge programs to correctly identify what info should go where and what can be left out, again with varying degrees of success.[4] For these reasons, OCR vendors and IDPs tend to specialize in different geographic areas and different document types. Throw one a European passport and you may be all set, but throw the same program a utility bill in Arabic and data capture accuracy may drop to unusable levels. You need to be able to identify the right solution for your needs, which likely means different solutions for different products and customers.

| Extracted Characters | | | | | |
|---|---|---|---|---|---|
| **Image Category** | **Existing Characters** | **Google Docs OCR** | **Tesseract** | **ABBY FineReader** | **Transym** |
| Digital Images | 1,834 | 1,613 (87.95%) | 1,539 (83.91%) | 1,528 (83.31%) | 1,463 (79.77%) |
| Machine-written characters | 703 | 569 (80.94%) | 549 (78.09%) | 574 (81.65%) | 554 (78.81%) |
| Hand-written characters | 2,036 | 1,254 (61.59%) | 984 (48.33%) | 1,204 (59.14%) | 960 (47.15%) |
| Black and white images | 71 | 69 (97.19%) | 69 (97.19%) | 65 (91.55%) | 61 (85.92%) |
| Multi-oriented text strings | 106 | 68 (64.15%) | 30 (28.3%) | 75 (70.75%) | 23 (21.7%) |
| PDF files | 15,693 | 15,409 (98.19%) | 14,121 (89.98%) | 15,376 (97.98%) | 14,133 (90%) |
| Multilingual text images | 3,597 | 2,831 (78.7%) | 2,474 (68.78%) | 2,799 (77.81%) | 1,740 (48.37%) |

## Protection plus precision: AI as the first step in document intake

Resistant AI is different: we check the documents you're presented with in a content-agnostic way. We **don't** focus on the content—that's the realm of OCR and IDPs. We focus instead on the image or document itself, how it's been created, and the way it's been presented to your platform. This crucial first step verifies whether a document is usable at all down the line while flagging manipulation and potential fraud attempts at the source. So instead of a specialist in a region or document type, you get a specialist in fraud.

Simultaneously, our quick-to-create models can be the key to a more flexible document processing workflow that directs a document to the OCR or IDP solution best suited to a specific task. This can help you get the most accurate data from each document while preventing wasted costs incurred by sending documents to ill-suited processors.

The conclusion is simple: many IDP options, one fraud solution.

[2] Cem Dilmegani, "Best OCR: Benchmark on Text Extraction / Capture Accuracy" (AIMultiple.com, 2022)
[3] Tafti, et al., "OCR as a Service: An Experimental Evaluation of Google Docs OCR, Tesseract, ABBYY FineReader, and Transym" (International Symposium on Visual Computing, 2016)
[4] Mariusz Kossakowski, "IDP Solutions - Custom Development vs Ready-To-Use Software" (Netguru, 2022)

## Get in touch and request a demo now